
Domain Adaptation by Mixture of Alignments of Second- or Higher-Order Scatter Tensors

Piotr Koniusz^{*,1,2}Yusuf Tas^{*,1,2}Fatih Porikli^{1,2}

¹Data61/CSIRO, ²Australian National University
 firstname.lastname@{data61.csiro.au, anu.edu.au}

Abstract

In this paper, we propose an approach to the domain adaptation, dubbed Second- or Higher-order Transfer of Knowledge (So-HoT), based on the mixture of alignments of second- or higher-order scatter statistics between the source and target domains. The human ability to learn from few labeled samples is a recurring motivation in the literature for domain adaptation. Towards this end, we investigate the supervised target scenario for which few labeled target training samples per category exist. Specifically, we utilize two CNN streams: the source and target networks fused at the classifier level. Features from the fully connected layers fc7 of each network are used to compute second- or even higher-order scatter tensors; one per network stream per class. As the source and target distributions are somewhat different despite being related, we align the scatters of the two network streams of the same class (within-class scatters) to a desired degree with our bespoke loss while maintaining good separation of the between-class scatters. We train the entire network in end-to-end fashion. We provide evaluations on the standard Office benchmark (visual domains), RGB-D combined with Caltech256 (depth-to-rgb transfer) and Pascal VOC2007 combined with the TU Berlin dataset (image-to-sketch transfer). We attain state-of-the-art results.

1 Introduction

Domain adaptation and transfer learning are the problems widely studied in computer vision and machine learning communities [1, 26]. They are directly inspired by the human cognitive abilities of generalizing to new concepts from very few data samples (cf. training from scratch on over a million of labeled images of the ImageNet dataset [29]). From psychological point of view, transfer of learning is “*the dependency of human conduct, learning or performance on prior experience*”. This problem was introduced in 1901 under a notion of “*transfer of particle*” [44]. In machine learning, transfer learning (or inductive learning) concerns “*storing knowledge gained while solving one problem and applying it to a different but related problem*” [43]. In practical computer vision and machine learning systems, transfer learning refers to “*an ability of a system to recognize and apply knowledge and skills learned in previous tasks to novel tasks or new domains, which share some commonality*”. In general, given a new (target) task, the arising question is how to identify the commonality between this task and previous (source) tasks, and transfer knowledge from the previous tasks to the target one. Therefore, one has to address three questions: what to transfer, how to transfer, and when to transfer [36].

In what follows, we propose an approach to the domain adaptation, dubbed Second- or Higher-order Transfer of Knowledge (*So-HoT*), based on the mixture of alignments of second- or higher-order scatter statistics between the source and target domains. Specifically, we utilize second- or higher-order scatter tensors, one per each network stream per class, such that the first stream corresponds to the source domain while the second to the target. The scatters are built from the feature vectors

*Both authors contributed equally.

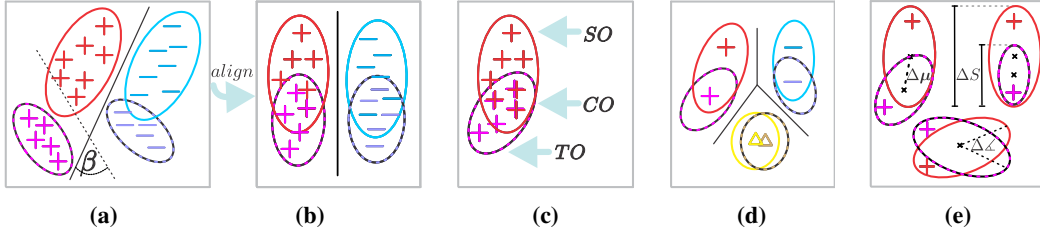


Figure 1: Our alignment problem. In Figure 1a, a two-class toy problem with positive and negative samples (+) and (−) is given. The solid and dashed ellipses indicate the source and target domain distributions. The two hyperplane lines that separate (+) from (−) on the target data indicate large uncertainty (denoted as β) in the optimal orientation for the target problem. Figure 1b shows that the source and target distributions can be aligned enough to separate well two classes for both the source and target problems. Figure 1c shows that partially aligned distributions have the commonality (CO) as well as the source and target specific parts (SO) and (TO) that represent dissimilarity between the source and target. Figure 1d depicts a multi-class problem. Beside of partially aligned means, the orientations of the source and target distributions are allowed to partially differ – as a result, they *i.e.* fit better into the piece-wise linear decision boundary. Figure 1e shows that differences in means $\Delta\mu$, scale/shear ΔS and orientation $\Delta\Delta$ of within-class scatters are all part of the alignment process.

produced by the *fc7* layer of AlexNet [21]. We propose that, as the source and target distributions are only partially related by their commonality, the scatters of the same class from both streams (*within-class scatters*) should be aligned to a desired degree to capture this commonality as an overlap between parts of the two distributions. At the same time, to achieve high classification accuracy, we maintain good separation between the scatters representing different classes (*between-class scatters*). We devise a simple loss that brings each pair of within-class scatters closer in terms of their covariances as well as their corresponding means. Therefore, the CNN parameters stored by convolutional filters and weights of the target network regularized by the source data in this end-to-end fashion must produce statistics consistent with the source network. We view such a regularization paradigm as being motivated by the theory of privileged learning [39]. In our case, the statistics of the source network regularize the target (and vice-versa) whilst in the privileged learning, the side information regularizes the solution dictated by the empirical loss evaluated on the main data samples. See Figures 1 and 2 for illustrative examples.

Furthermore, as distributions of the source and target domains may require different level of alignment per class (the commonality depends on the class label), we investigate not only an unweighted alignment loss (class-independent level of alignment) but also its weighted counterpart which learns one weight per class (class-specific levels of alignment).

Additionally, as we work with second- or higher-order tensors, we propose a kernelized variant of our alignment loss which provides computational speed-ups for typical domain adaptation datasets.

To summarize, our main contributions are: i) a novel loss that we call *So-HoT*, which defines the commonality between the source and target domains as the mixture of alignments of second- or higher-order scatter tensors, ii) unweighted and weighted variants of alignments, and iii) a fast kernelized alternative of our alignment loss.

Next, we detail the notion of domain adaptation and transfer learning, review the related literature and explain how our work differs from the state-of-the-art approaches.

2 Related Work

Domain adaptation assumes that the transfer of knowledge takes place among two or more distinct domains *e.g.*, e-commerce reviews and biomedical data. In contrast, transfer learning utilizes the same domain *e.g.*, images of natural scenes with related but different distributions where the goal may be to learn objects of a new class while leveraging other already learned classes [35, 36]. Not surprisingly, these both notions are often interchangeable *e.g.*, natural images and sketches have related distributions but they come from distinct domains at the same time. Another example is a so-called domain shift *e.g.*, bicycle in natural images vs. on-line retailer galleries. Transfer of knowledge may vary from simply carrying over discriminative information from a source to target domain under the same set of classes to inferring a solution to a new distinct task from a set of former ones [35, 16]. Domain adaptation comes in many flavors. Single- or multiple-source [4] setups are

possible *e.g.*, single stream of natural images vs. multiple streams supplied with photos of objects: on cluttered backgrounds, on a clear background, in a daytime or night setting, or even in multi-spectral setting. Moreover, the problem in hand may be homogeneous or heterogeneous [36, 45] in nature *e.g.*, identical source and target representations using RGB images vs. a source represented by a CNN trained on images [9, 21] and a target using an RNN-inspired [15] LSTM [14] which is trained on text corpora. The architecture in use may be shallow [5, 34] or deep [10] such that the commonality is established only at the classifier level or across entire source and target networks, respectively. Noteworthy is also recent trend in the CNN fine-tuning which by itself is a powerful domain adaptation and transfer learning tool [12, 31] which requires large training datasets. Moreover, domain adaptation and transfer learning address problems such as: learning new categories from few annotated samples (supervised domain adaptation [3, 38]), utilizing available unlabeled data (unsupervised [34, 10] or semi-supervised domain adaptation [5, 38]), recognizing new categories in embedded spaces *e.g.*, attribute-based, without any training samples (zero-shot learning [23]).

In this paper, we investigate the case of a deep supervised single-source domain adaptation which can be easily extended to the multi-source and heterogeneous cases.

The Commonality. Deep learning has been used in the context of domain adaptation in recent works *e.g.*, [38, 10, 3, 42, 22, 37], just to name a few. These works differ in how they establish the so-called commonality between domains. In [38], the authors propose to align both domains via the cross entropy which “maximally confuses” both domains for supervised and semi-supervised settings. In [10], an unsupervised approach utilizes the assumption that predictions must be made based on features which cannot discriminate between the source and target domains. Specifically, they minimize a trade-off between the so-called source risk and the empirical divergence to find examples in the source domain indistinguishable from the target samples.

Our work differs from these approaches in that we define the commonality as the desired degree of overlap between the second- or higher-order scatters of the source and target. After such an alignment, we allow the non-overlapping tails of distributions to also guide the learning process to learn a more general classifier (*i.e.* avoid the domain-specific bias).

Moreover, in [3], the authors capture the “interpolating path” between the source and target domains using linear projections into a low-dimensional subspace which lies on the Grassman manifold. In [42], the authors propose to learn the transformation between the source and target by the deep model regression network. These two approaches assume that the source representation can be interpolated or regressed into the target as, given the nature of CNNs, they can approximate highly non-linear functions.

Our model differs in that our source and target network streams co-regularize each other to produce the commonality between the source and target distributions and accommodate the domain-specific parts that should not be aligned.

For visual domains, the commonality can be captured in the spatially-local sense. In [37], the authors utilize so-called “domainness maps” which capture locally the degree of domain specificity. Similarly, in [22], the authors extract local patches of varying sizes at process each of these patches via CNNs. Our work is orthogonal to these techniques. We represent the commonality globally, however, our ideas could also be applied in a spatially-local setting.

Correlation Methods. Some recent works use correlation between the source and target distributions. Inspired by the Canonical Correlation Analysis (CCA), the authors of [45] utilize a correlation subspace as a joint representation for associating the data across different domains. They also use kernelized CCA. In [34], the authors propose an unsupervised domain adaptation by the correlation alignment.

Our work is somewhat related in that we utilize second-order statistics. However, we perform a partial alignment of class-specific source and target distributions to define the commonality (partial intersection of scatters) in the supervised setting. We also align partially the distribution means while the above unsupervised approaches use zero-centered feature vectors and the full alignment of the generic (c.f. class-specific) source and target distributions. We propose how to learn the degree of alignment in an end-to-end fashion and introduce the kernelized loss between the second- or higher-order scatter tensors; all being novel propositions.

Tensor Methods. Correlation approaches outlined above use second-order scatter matrices which are tensors of order $r=2$. In this work, we also investigate the applicability of higher-order scatters

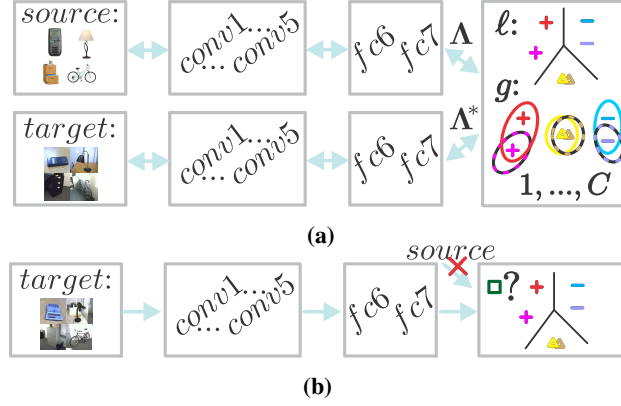


Figure 2: The pipeline. Figure 2a shows the source and target network streams which merge at the classifier level. The classification and alignment loss ℓ and g take the data Λ and Λ^* from both streams and participate in end-to-end learning. At the test time, we use the target stream and the trained classifier as in Figure 2b.

$r \geq 3$ for alignment. Third-order tensors have been found useful for various vision tasks. For example, spatio-temporal third-order tensor on video data is proposed for action analysis in [17], non-negative tensor factorization is used for image denoising in [32], tensor textures are proposed for texture rendering in [41], and higher order tensors are used for face recognition in [40]. A survey of multi-linear algebraic methods for tensor subspace learning is available in [27]. The above applications use a single tensor, while our goal is to use tensors as the domain- and class-specific representations, similar to the sum-kernel approaches [20, 18, 19], and apply them to alignment tasks.

3 Background

In this section, we review our notations and the necessary background on scatter tensors, polynomial kernels and their linearizations, which will be useful in deriving our mixture of alignments of second- or higher-order scatter tensors.

3.1 Notations

Let $\mathbf{x} \in \mathbb{R}^d$ be a d -dimensional feature vector. Then, we use $\mathcal{X} = \uparrow \otimes_r \mathbf{x}$ to denote the r -mode super-symmetric rank-one tensor \mathcal{X} generated by the r -th order outer-product of \mathbf{x} , where the element of $\mathcal{X} \in \mathfrak{S}_{\times r}^d$ at the (i_1, i_2, \dots, i_r) -th index is given by $\prod_{j=1}^r x_{i_j}$. \mathcal{I}_N stands for the index set $\{1, 2, \dots, N\}$. We denote the space of super-symmetric tensors of dimension $d \times \dots \times d$ with r modes as $\mathfrak{S}_{\times r}^d \subset \mathbb{R}^{\times_r d}$, where $\mathbb{R}^{\times_r d}$ is the space of tensors $\mathbb{R}^{d \times \dots \times d}$ with r modes. The Frobenius norm of tensor is given by $\|\mathcal{X}\|_F = \sqrt{\sum_{i_1, i_2, \dots, i_r} \mathcal{X}_{i_1, i_2, \dots, i_r}^2}$, where $\mathcal{X}_{i_1, i_2, \dots, i_r}$ represents the (i_1, i_2, \dots, i_r) -th element of \mathcal{X} . Similarly, the inner-product between two tensors \mathcal{X} and \mathcal{Y} is given by $\langle \mathcal{X}, \mathcal{Y} \rangle = \sum_{i_1, i_2, \dots, i_r} \mathcal{X}_{i_1, i_2, \dots, i_r} \cdot \mathcal{Y}_{i_1, i_2, \dots, i_r}$. Using Matlab style notation, the (i_3, \dots, i_r) -th slice of \mathcal{X} is given by $\mathcal{X}_{:, :, i_3, \dots, i_r}$. The space of positive semi-definite matrices is \mathcal{S}_+^d . Lastly, $\mathbf{1}$ denotes a vector with all coefficients equal one.

3.2 Second- or Higher-order Scatter Tensors

We define a scatter tensor of order r as a mean-centered TOSST representation [18]:

Definition 1. Suppose $\phi_n \in \mathbb{R}^d, \forall n \in \mathcal{I}_N$ represent some data vectors, then a scatter tensor $\mathcal{X} \in \mathfrak{S}_{\times r}^d$ of order r on these data vectors is given by:

$$\mathcal{X} = \frac{1}{N} \sum_{n=1}^N \uparrow \otimes_r (\phi_n - \mu) \quad \text{and} \quad \mu = \frac{1}{N} \sum_{n=1}^N \phi_n. \quad (1)$$

In our supervised domain adaptation setting, the scatter tensors are obtained via applying (1) on the class-specific data vectors such as outputs of the *fc7* layer of AlexNet.

The following properties of the scatter tensors are worth noting (see [18] for proofs):

Proposition 1. For a scatter tensor $\mathcal{X} \in \mathbb{S}_{\times r}^d$, we have:

1. *Super-Symmetry:* $\mathcal{X}_{i_1, i_2, \dots, i_r} = \mathcal{X}_{\Pi(i_1, i_2, \dots, i_r)}$ for indexes (i_1, i_2, \dots, i_r) and their any permutation Π . The number of unique coefficients of \mathcal{X} is $\binom{d+r-1}{r}$.
2. *Every slice is at least positive semi-definite* for any even order $r \geq 2$ and $\mathcal{X}_{::, i_3, \dots, i_r} \in \mathcal{S}_+^d, \forall (i_3, \dots, i_r) \in \mathcal{I}_d$. For $r=2$, tensor \mathcal{X} also is a covariance matrix.
3. *Indefiniteness* for any odd order $r \geq 1$, i.e., under a CP decomposition [25], it can have positive, negative, or zero entries in its core-tensor.

Due to the indefiniteness of tensors of odd orders and potential rank deficiency, we restrict ourselves to work with the Euclidean distance between such scatter representations. Also, as the number of unique coefficients of \mathcal{X} is of order $\sim d^r$, which is prohibitive for $r \geq 3$, we propose a light-weight kernelized variant of the Euclidean distance which avoids explicit use of tensors. The following easily verifiable two results will come handy in the sequel:

Proposition 2. Suppose we want to evaluate the Frobenius norm between tensors $\mathcal{X}, \mathcal{X}^* \in \mathbb{S}_{\times r}^d$, then it holds that:

$$\|\mathcal{X} - \mathcal{X}^*\|_F^2 = \langle \mathcal{X}, \mathcal{X} \rangle - 2 \langle \mathcal{X}, \mathcal{X}^* \rangle + \langle \mathcal{X}^*, \mathcal{X}^* \rangle. \quad (2)$$

Proof. \mathcal{X} and \mathcal{X}^* can be vectorized and the Frobenius norm replaced by the ℓ_2 -norm for which the above expansion is known to hold. Then, folding back vectors and re-applying the Frobenius norm in place of the ℓ_2 -norm completes the proof. \square

Proposition 3. Suppose $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ are two arbitrary vectors, then for an ordinal $r > 0$, we have:

$$\langle \mathbf{x}, \mathbf{y} \rangle^r = \langle \uparrow \otimes_r \mathbf{x}, \uparrow \otimes_r \mathbf{y} \rangle. \quad (3)$$

Moreover, for sets of vectors $\mathbf{x}_n, \mathbf{y}_{n'} \in \mathbb{R}^d$, we have:

$$\sum_n \sum_{n'} \langle \mathbf{x}_n, \mathbf{y}_{n'} \rangle^r = \left\langle \sum_n \uparrow \otimes_r \mathbf{x}_n, \sum_{n'} \uparrow \otimes_r \mathbf{y}_{n'} \right\rangle. \quad (4)$$

Proof. The expansion in (3) is derived in [20] while (4) can be verified due to bilinear properties of the dot-product. \square

4 Proposed Approach

In this section, we first formulate the problem of mixture of alignments of second- or higher-order scatter tensors, which precedes an exposition to our next two contributions, that is, a weighted mixture of alignments and a kernelized approach which avoids explicit evaluations of scatters.

4.1 Problem Formulation

Suppose \mathcal{I}_N and \mathcal{I}_{N^*} are the indexes of N source and N^* target training data points. \mathcal{I}_{N_c} and $\mathcal{I}_{N_c^*}$ are the class-specific indexes for $c \in \mathcal{I}_C$, where C is the number of classes. Suppose we have feature vectors from *fc7* in the source network stream, one per image, and associated with them labels. Such pairs are given by $\mathbf{A} \equiv \{(\phi_n, y_n)\}_{n \in \mathcal{I}_N}$, where $\phi_n \in \mathbb{R}^d$ and $y_n \in \mathcal{I}_C, \forall n \in \mathcal{I}_N$, as shown in Figure 2a. For the target data, by analogy, we define pairs $\mathbf{A}^* \equiv \{(\phi_n^*, y_n^*)\}_{n \in \mathcal{I}_{N^*}}$, where $\phi_n^* \in \mathbb{R}^d$ and $y_n^* \in \mathcal{I}_C, \forall n \in \mathcal{I}_{N^*}$. Class-specific sets of feature vectors are given as $\Phi_c \equiv \{\phi_n^c\}_{n \in \mathcal{I}_{N_c}}$ and $\Phi_c^* \equiv \{\phi_n^{c*}\}_{n \in \mathcal{I}_{N_c^*}}, \forall c \in \mathcal{I}_C$. Then, $\Phi \equiv (\Phi_1, \dots, \Phi_C)$ and $\Phi^* \equiv (\Phi_1^*, \dots, \Phi_C^*)$. Note that we use the asterisk symbol written in superscript (e.g. ϕ^*) to denote variables associated with the target network whilst the source-related and generic variables have no such indicator. Below, we formulate our problem as a trade-off between the classifier loss ℓ and the alignment loss g which acts on the scatter tensors and related to them means:

$$\begin{aligned} \arg \min_{\mathbf{W}, \mathbf{b}, \Theta, \Theta^*} \quad & \ell(\mathbf{W}, \mathbf{b}, \mathbf{A} \cup \mathbf{A}^*) + \lambda \|\mathbf{W}\|_F^2 \\ \text{s. t. } \quad & \|\phi_n\|_2^2 \leq \tau, \\ & \|\phi_{n'}^*\|_2^2 \leq \tau, \\ & \forall n \in \mathcal{I}_N, n' \in \mathcal{I}_{N^*} \end{aligned} \quad + \underbrace{\frac{\sigma_1}{C} \sum_{c \in \mathcal{I}_C} \|\mathcal{X}_c - \mathcal{X}_c^*\|_F^2 + \frac{\sigma_2}{C} \sum_{c \in \mathcal{I}_C} \|\mu_c - \mu_c^*\|_2^2}_{g(\Phi, \Phi^*)}. \quad (5)$$

For ℓ , we use a generic loss used by CNNs, say Softmax. The matrix $\mathbf{W} \in \mathbb{R}^{d \times C}$ contains unnormalized probabilities (c.f. hyperplane of SVM), $\mathbf{b} \in \mathbb{R}^C$ is the bias term, and λ is the regularization constant. Moreover, the union $\mathbf{A} \cup \mathbf{A}^*$ of the source and target training data reveals that we train one universal classifier for both domains¹. In Equation (5), separating the class-specific distributions is addressed by ℓ while bringing closer the within-class scatters of both network streams is handled by g (as Figure 2 shows). Specifically, our loss g depends on two sets of variables $(\mathcal{X}_1(\Phi_1), \dots, \mathcal{X}_C(\Phi_C)), (\mu_1(\Phi_1), \dots, \mu_C(\Phi_C))$ and $(\mathcal{X}_1^*(\Phi_1^*), \dots, \mathcal{X}_C^*(\Phi_C^*)), (\mu_1^*(\Phi_1^*), \dots, \mu_C^*(\Phi_C^*))$ – one set per network stream. Feature vectors $\Phi(\Theta)$ and $\Phi^*(\Theta^*)$ depend on the parameters of the source and target network streams Θ and Θ^* that we optimize over *e.g.*, they represent coefficients of convolutional filters and weights of *fc* layers. $\mathcal{X}_c, \mathcal{X}_c^*, \mu_c$ and μ_c^* denote the scatter tensors and means, respectively, one tensor/mean pair per network stream per class, evaluated as in (1). Lastly, σ_1 and σ_2 control the overall degree of the scatter and mean alignment, τ constraints the ℓ_2 -norm of feature vectors (needed if λ is low). Derivatives of loss g are given in Appendix B.

In this work, we make an educated assumption that highly non-linear CNN streams are able to rotate the within-class scatters sufficiently, as our loss dictates, so that the commonality becomes represented as a partial overlap of two scatters. Such an assumption is common in *i.e.* [3, 42].

4.2 Weighted Alignment Loss

Below we propose a weighted variant of alignment loss g which incorporates class-specific weights $\zeta \in \mathbb{R}^C$ and $\bar{\zeta} \in \mathbb{R}^C$ that adjust the degree of alignment per class between the within-class scatters as well as related to them means, respectively. To achieve this objective, we formulate the alternative loss g as follows:

$$g(\Phi, \Phi^*, \zeta, \bar{\zeta}) = \frac{\sigma_1}{C} \sum_{c \in \mathcal{I}_C} \zeta_c \|\mathcal{X}_c - \mathcal{X}_c^*\|_F^2 + \frac{\sigma_2}{C} \sum_{c \in \mathcal{I}_C} \bar{\zeta}_c \|\mu_c - \mu_c^*\|_2^2 + \alpha_1 \|\zeta - \mathbf{1}\|_2^2 + \alpha_2 \|\bar{\zeta} - \mathbf{1}\|_2^2, \quad (6)$$

where α_1 and α_2 control the degree of weight deviation. When using the weighted alignment, we replace the corresponding loss in Eq. (5) by the alignment loss g defined in (6). Then, we additionally minimize (5) over ζ and $\bar{\zeta}$.

4.3 Kernelized Alignment Loss

Evaluating scatter tensors at each iteration of the gradient descent can be costly, even if using co-variances ($r = 2$), as the typical size of feature vectors of *fc7* is $d = 4096$. Below we propose an efficient kernelization of the Frobenius norm.

Proposition 4. *The inner-product of scatter tensors $\mathcal{X}, \mathcal{Y} \in \mathbb{S}_{\times^r}^d$, which are defined in (1), can be expressed implicitly as a sum of entries of a polynomial kernel $\bar{\mathbf{K}}^r \in \mathbb{R}^{N \times N^*}$, where $\bar{\mathbf{K}}_{nn'}^r = \langle \mathbf{x}_n - \mu, \mathbf{y}_{n'} - \mu^* \rangle^r$, and $\mathbf{x}_n \in \mathbb{R}^d, \forall n \in \mathcal{I}_N$ and $\mathbf{y}_{n'} \in \mathbb{R}^d, \forall n' \in \mathcal{I}_{N^*}$ are some N and N^* feature vectors (used to build \mathcal{X} and \mathcal{Y}), μ and μ^* are their means. Then:*

$$\langle \mathcal{X}, \mathcal{Y} \rangle = \frac{1}{NN^*} \sum_n \sum_{n'} \langle \mathbf{x}_n - \mu, \mathbf{y}_{n'} - \mu^* \rangle^r = \frac{1}{NN^*} \mathbf{1}^T \bar{\mathbf{K}}^r \mathbf{1}. \quad (7)$$

Proof. Substituting $\mathbf{x}_n - \mu$ and $\mathbf{y}_{n'} - \mu^*$ into Proposition 3, the proof follows. \square

Proposition 5. *Suppose we have polynomial kernels $\mathbf{K}^r \in \mathbb{R}^{N \times N}$, $\bar{\mathbf{K}}^r \in \mathbb{R}^{N^* \times N^*}$ and $\bar{\bar{\mathbf{K}}}^r \in \mathbb{R}^{N \times N^*}$ defined as $K_{nn'}^r = \langle \mathbf{x}_n - \mu, \mathbf{x}_{n'} - \mu \rangle^r$, $\bar{K}_{nn'}^r = \langle \mathbf{y}_{n'} - \mu^*, \mathbf{y}_{n'} - \mu^* \rangle^r$ and $\bar{\bar{K}}_{nn'}^r = \langle \mathbf{x}_n - \mu, \mathbf{y}_{n'} - \mu^* \rangle^r$, where $\mathbf{x}_n, \mathbf{y}_{n'}, \mu, \mu^*, N, N^*$ are defined as in Proposition 4. The Frobenius norm between two scatter tensors $\mathcal{X}, \mathcal{Y} \in \mathbb{S}_{\times^r}^d$, which are defined in (1), can be expressed implicitly as:*

$$\|\mathcal{X} - \mathcal{X}^*\|_F^2 = \frac{1}{N^2} \mathbf{1}^T \mathbf{K}^r \mathbf{1} + \frac{1}{N^{*2}} \mathbf{1}^T \bar{\mathbf{K}}^r \mathbf{1} - \frac{2}{NN^*} \mathbf{1}^T \bar{\bar{\mathbf{K}}}^r \mathbf{1}. \quad (8)$$

Proof. Combining Proposition 2 with 4, the proof follows. \square

¹ In the heterogeneous setting, we use two domain-specific classifiers *i.e.*, $\ell(\mathbf{W}, \mathbf{b}, \mathbf{A}) + \ell(\mathbf{W}^*, \mathbf{b}^*, \mathbf{A}^*) + \lambda \|\mathbf{W}\|_F^2 + \lambda^* \|\mathbf{W}^*\|_F^2 + \beta' \|\mathbf{W} - \mathbf{W}^*\|_F^2$. Otherwise, we avoid this method as it may overfit more.

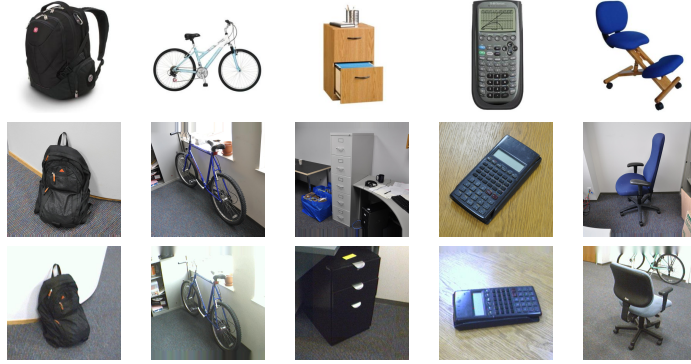


Figure 3: The Office dataset. Top, middle and bottom rows show examples from the Amazon, DSLR, and Webcam domains.

Derivatives of (8) are in Appendix C. Equation (8) can be evaluated on class-specific feature vectors and substituted directly into the loss functions in (5) and (6). This way, we obtain two different regimes for evaluating the Frobenius norm on the scatter tensors: one explicit and one kernelized; both exhibiting different strengths as detailed below.

Complexity. The Frobenius norm on the scatter tensors has complexity $\mathcal{O}((N+N^*+1)D)$, where $D = \binom{d+r-1}{r}$ as detailed in Proposition 1. The kernelized variant proposed above has complexity $\mathcal{O}((N^2 + NN^* + N^*N^*)(d+\rho))$, where $\rho \leq \log r$ estimates the complexity of “rising x to the power of r ”. As $\rho \ll d$, its cost is negligible and can be safely left out from the above analysis.

It is easy to verify that, for the standard domain adaptation problems with $N=20$ source and $N^*=3$ target training points per class, $d=4096$ and $r=2$, explicit evaluations of the Frobenius norm are $\sim 52\times$ slower than the proposed by us kernelized substitute. For the same scenario but with the scatter tensor of order $r=3$, explicit evaluations of the Frobenius norm are not tractable, as they take $\sim 143000\times$ more time than the kernelized substitute, which demonstrates the clear benefit of our approach.

5 Experiments

In this section, we present experiments demonstrating the usefulness of our framework. We start by describing datasets we use in evaluations.

5.1 Datasets

Office dataset. A popular dataset for evaluating algorithms against the effect of domain shift is the Office dataset [30] which contains 31 object categories in three domains: Amazon, DSLR and Webcam. The 31 categories in the dataset consist of objects commonly encountered in office settings, such as keyboards, file cabinets, and laptops. The Amazon domain contains on average 90 images per class and 2817 images in total. As these images were captured from a website of online merchants, they are captured against clean background and at a unified scale. The DSLR domain contains 498 low-noise high resolution images (4288×2848). There are 5 objects per category. Each

		$\mathcal{A} \rightarrow \mathcal{W}$	$\mathcal{A} \rightarrow \mathcal{D}$	$\mathcal{W} \rightarrow \mathcal{A}$	$\mathcal{W} \rightarrow \mathcal{D}$	$\mathcal{D} \rightarrow \mathcal{A}$	$\mathcal{D} \rightarrow \mathcal{W}$	acc.
DLID	[3]	51.9	-	-	89.9	-	78.2	73.33
DeCAF ₆ S+T	[6]	80.7 \pm 2.3	-	-	-	-	94.8 \pm 1.2	87.75
DaNN	[11]	53.6 \pm 0.2	-	-	83.5 \pm 0.0	-	71.2 \pm 0.0	69.43
Source CNN	[38]	56.5 \pm 0.3	64.6 \pm 0.4	42.7 \pm 0.1	93.6 \pm 0.2	47.6 \pm 0.1	92.4 \pm 0.3	66.23
Target CNN	[38]	80.5 \pm 0.5	81.8 \pm 1.0	59.9 \pm 0.3	81.8 \pm 1.0	59.9 \pm 0.3	80.5 \pm 0.5	74.06
Source+Target CNN	[38]	82.5 \pm 0.9	85.2 \pm 1.1	65.2 \pm 0.7	96.3 \pm 0.5	65.8 \pm 0.5	93.9 \pm 0.5	81.48
Dom. Conf.+Soft Labs.	[38]	82.7 \pm 0.8	86.1 \pm 1.2	65.0 \pm 0.5	97.6\pm0.2	66.2 \pm 0.3	95.7\pm0.5	82.22
Source+Target CNN (S+T)		82.4 \pm 2.0	85.5 \pm 0.9	65.1 \pm 1.4	95.8 \pm 0.8	66.0 \pm 1.2	94.3 \pm 0.6	81.53
Second-order (S_o)		84.5\pm1.7	86.3\pm0.8	65.7\pm1.7	97.5 \pm 0.7	66.5\pm1.0	95.5 \pm 0.6	82.68

Table 1: Comparison of our second-order alignment loss (S_o) to the state of the art on the Office dataset.

object was captured from different viewpoints on average 3 times. For Webcam, the 795 images of low resolution (640×480) exhibit significant noise and color as well as white balance artifacts. Otherwise, 5 objects per category were also used in the capturing process. Figure 3 illustrates the three domains. We distinguish the following six domain shifts: Amazon-Webcam ($\mathcal{A} \rightarrow \mathcal{W}$), Amazon-DSLR ($\mathcal{A} \rightarrow \mathcal{D}$), Webcam-Amazon ($\mathcal{W} \rightarrow \mathcal{A}$), Webcam-DSLR ($\mathcal{W} \rightarrow \mathcal{D}$), DSLR-Amazon ($\mathcal{D} \rightarrow \mathcal{A}$) and DSLR-Webcam ($\mathcal{D} \rightarrow \mathcal{W}$).

Unless stated otherwise, for each of the above domain shifts, we evaluate across 10 randomly chosen data splits. We follow the standard protocol for this dataset and, for each training source split, we sample 20 images per category for the Amazon domain and 8 examples per category for the DSLR and Webcam domains. From the training target splits, we sample 3 images per class per split per domain. We present results for the supervised setting and report accuracies on the remaining target images, as the standard protocol for this dataset suggests.

RGB-D-Caltech256 dataset. The RGB-D [24] and Caltech256 [13] datasets have been used as the source and target for evaluations of unsupervised domain adaptation problems [2, 28]. We use the 10 classes that are common between the two datasets *e.g.*, calculator, cereal box, coffee mug, ball, tomato. We use 50 and 5/10 images per class in the source and target domains for the supervised setting. We test on the remaining target samples. We report the mean average accuracy over 5 data splits, that is, we select randomly the source and target data samples for each split.

Pascal VOC2007-TU Berlin dataset. Transfer from Pascal VOC2007 [8] to TU Berlin [7] (images-to-sketches transfer) has never been attempted yet in domain adaptation to our best knowledge. We utilize 50 and 3 source and target training samples per class, respectively, and the 14 classes that are common between the source and target datasets. We perform testing on the remaining target data. We report the mean average accuracy over 5 data splits.

5.2 Experimental Setup

In each stream, we employ the AlexNet architecture [21] which was pre-trained on the ImageNet dataset [29] for the best results. At the training and testing time, we use the pipelines shown in Figures 2a and 2b, respectively. Where stated, we use the 16-layer VGG model [33] per stream to quantify the impact of different CNN models on our algorithm. We set non-zero learning rates on the fully-connected and the last two convolutional layers of the two streams.

On the RGB-D-Caltech256 dataset, we use the RGB images from Caltech256 as the target domain. In contrast to [2, 28] which use both the RGB data and depth maps as a source, we adapt our source stream based on AlexNet to use only the depth data from the RGB-D dataset – this helps us isolate performance of our algorithm in case of distinct heterogeneous domains. As these both domains are very different from each other, we apply two classifiers – one per network stream (see the footnote¹ in Section 4).

We evaluate three variants of our Second- or Higher-order Transfer of Knowledge (*So-HoT*) approach: unweighted and weighted second-order alignment losses (*So*) and (*So+ ζ*) for $r=2$, and the third-order loss (*To*) for $r=3$. The model parameters were selected by cross-validation.

5.3 Comparison to the State of the Art

We apply our algorithm on the Office dataset. Table 1 presents results for the six domain shifts. Our second-order alignment loss (*So*) is compared against the baseline (*S+T*) for which the source and target training samples were used together to fine-tune a standard CNN network. As can be seen, our method outperforms such a baseline as well as recent approaches such as *Domain Confusion with Soft Labels* and fine-tuning on the source or target data. Some approaches in the literature *e.g.*, work in [28], report marginally higher results *e.g.* 86.1% on $\mathcal{A} \rightarrow \mathcal{W}$ but they use all the source and half of the target data for training. In this setting, Table 2 shows that our second-order approach (*So*) outperforms (*S+T*) by 1% and unsupervised [28] by 10.5%.

	sp1	sp2	sp3	sp4	sp5	average acc.
<i>S+T</i>	96.2	98.1	93.1	96.2	95.0	95.72 ± 1.8
<i>So</i>	97.5	98.1	95.0	96.8	96.1	96.70 ± 1.2

Table 2: Second-order approach (*So*) vs. the combined source and target (*S+T*) trained on large source and target data in $\mathcal{A} \rightarrow \mathcal{W}$.

		sp1	sp2	sp3	sp4	sp5	sp6	sp7	sp8	sp9	sp10	avg. acc.
AlexNet	$S+T$	85.9	79.3	80.6	84.2	83.0	82.9	81.3	82.0	84.3	80.6	82.41 ± 2.0
	So	87.16	82.31	84.19	86.32	83.76	86.47	82.91	83.46	85.18	83.19	84.49 ± 1.6
	$So+\zeta$	87.74	82.90	84.47	86.75	84.76	86.47	83.20	83.90	85.47	83.48	84.91 ± 1.6
	To	87.3	83.5	83.9	86.0	84.1	86.3	82.6	83.6	86.04	85.2	84.85 ± 1.5
VGG	$S+T$	91.50	87.89	91.45	89.60	89.03	87.18	86.18	87.18	91.17	86.47	88.76 ± 1.9
	So	91.88	89.32	91.74	90.60	89.03	88.18	87.60	87.60	91.88	87.18	89.5 ± 1.8

Table 3: (Top) Evaluations with use of AlexNet streams. We compare our weighted alignment loss ($So+\zeta$) and third-order alignment loss (To) with our baseline fine-tuning on the combined source and target domains ($S+T$) and the second-order approach (So). (Bottom) Given the VGG streams, we compare our second-order approach (So) to the baseline ($S+T$). In all cases, we use the $\mathcal{A} \rightarrow \mathcal{W}$ domain shift, report accuracy (%) for all 10 splits and the average accuracy.

Weighted vs. Unweighted Alignment. In this experiment, we demonstrate the benefit of using the weighted alignment of the scatter matrices and their means on the $\mathcal{A} \rightarrow \mathcal{W}$ domain shift. Table 3 shows that our weighted second-order alignment loss ($So+\zeta$), introduced in (6), improves over our unweighted second-order alignment loss (So) from (5) by 0.4%. The increase is consistent on all 10 splits. As learning ζ and $\bar{\zeta}$ can be easily added into the code at no visible increase in computations, such a strategy appears beneficial.

In Figure 4, we show histograms of the ζ and $\bar{\zeta}$ weights over the 31 classes and the 10 splits. The histograms reveal that the levels of alignment of the scatter matrices and their means vary according to the Beta distribution. The means of these distributions are slightly below the desired mean value of one which indicates that, in this experiment, σ_1 and σ_2 from Equation (6) were initialized with values larger than needed. Also, their optimal values might have varied over time. Learning weights compensates for such effects.

Alignment of Third-order Scatter Tensors. Our kernelized loss in (8) admits alignment between higher-order scatter tensors which, beyond the scale/shear and orientation, capture higher-order statistical moments. In Table 3, we evaluate third-order alignment loss (To) which marginally outperforms our second-order approach (So) at no extra computational cost, as detailed in Section 4.3.

Performance on the VGG architecture. To evaluate effectiveness of our algorithm on other powerful networks, we follow the same pipeline as in Figure 2, except that we employ the pre-trained VGG [33] in place of AlexNet [21]. As VGG utilizes more parameters than AlexNet, we demonstrate in Table 3 (bottom) that applying our second-order alignment loss (So) on $\mathcal{A} \rightarrow \mathcal{W}$ improves performance compared to the baseline ($S+T$) even on such a powerful architecture. Without resorting to any data augmentations, we also outperform by 0.6% a multi-scale multi-patch CNN approach [22]. For completeness, we also evaluated our approach on the $\mathcal{A} \rightarrow \mathcal{D}$ domain shift. Our (So) scored $92.3 \pm 1.1\%$ which outperforms the baseline ($S+T$) of $89.8 \pm 1.4\%$ by 2.5%. Moreover, we also obtained $73.11 \pm 0.9\%$ accuracy on the $\mathcal{W} \rightarrow \mathcal{A}$ domain shift. In contrast, the baseline ($S+T$) scored $72.2 \pm 0.9\%$ accuracy. Despite we use a more powerful VGG architecture here (c.f. AlexNet) which benefits fine-tuning, we consistently outperform fine-tuning and obtain state of the art results on reported domains shifts.

Heterogeneous setting on RGB-D-Caltech256. On 5 data splits and 5 target samples, fine-tuning on the target domain (T) and the source combined with the target ($S+T$) yield 76.5 ± 1 and $76.8 \pm 1\%$ accuracy. Our second-order alignment loss (So) scores $77.84 \pm 1.2\%$. On 10 target samples, (T), ($S+T$) and (So) give 81.4 ± 1.3 , 81.5 ± 1 and $82.1 \pm 1.1\%$.

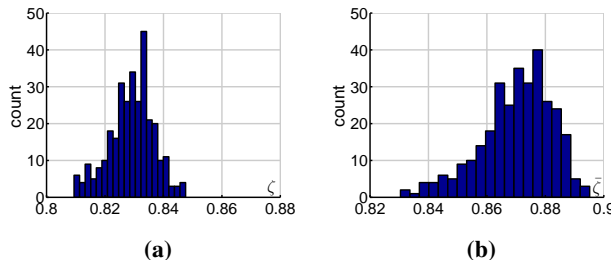


Figure 4: Histograms of the ζ and $\bar{\zeta}$ weights in plots 4a and 4b, learned on $\mathcal{A} \rightarrow \mathcal{W}$, show the level of alignment of the scatter matrices and their means according to the loss function in (6).

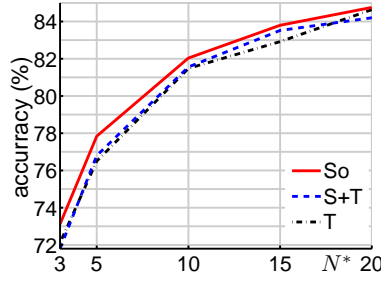


Figure 5: Our second-order algorithm (So) vs. the baseline fine-tuning on: i) the combined source and target domains ($S+T$) and ii) the target domain only (T). We use RGB-D-Caltech256 (heterogeneous setup). N^* is the number of target tr. samples per class.

	sp1	sp2	sp3	sp4	sp5	average acc.
$S+T$	58.86	63.43	63.14	59.14	68.71	62.66
T	59.86	63.43	64.14	57.86	67.0	62.46
So	60.57	63.28	64.28	59.14	69.71	63.40

Table 4: Pascal VOC2007-TU Berlin dataset. We use 5 splits and report accuracies on our method (So) vs. baselines ($S+T$) and (T).

Moreover, we verify the behavior of our second-order alignment loss (So) w.r.t. the varying number of target training samples. Figure 5 shows that the largest improvement of 1.24% over the baseline ($S+T$) is obtained for a small number $N^* = 3$. As N^* increases, the improvement over baselines becomes smaller. Such a trend is consistent with other works on domain adaptation [36]. In some cases, the baseline ($S+T$) performs worse than the fine-tuning on target only (T) which is known as so-called *negative transfer* [36]. For all $3 \leq N^* \leq 20$, our (So) outperforms baselines ($S+T$) and (T) which demonstrates robustness of our approach.

Heterogeneous setting on Pascal VOC2007-TU Berlin. Table 4 shows results on transfer from Pascal VOC2007 [8] to TU Berlin [7] (images-to-sketches transfer). As demonstrated in the table, our second-order approach (So) and the baselines ($S+T$) and (T) yield 63.4, 62.66 and 62.46%, respectively. Even in case of such a difficult heterogeneous problem, our algorithm still yields a consistent improvement.

6 Conclusions

We have presented an approach to domain adaptation by partial alignment of the within-class scatters to discover the commonality. The state-of-the-art results we obtain suggest that our simple strategy is effective despite challenges of domain adaptation. Moreover, the presented weighted approach and kernelized alignment loss improve the results and computational efficiency. Our method can be easily extended to multiple domains and other network architectures.

Acknowledgements. We thank Dr. Hongping Cai for our early discussions on the alignment loss, NVIDIA for the donation of GPUs and CSIRO for making the Bragg cluster available to us.

Appendices

A Sensitivity to parameters σ_1 and σ_2

Below we evaluate sensitivity of our second-order alignment loss (So) to the parameter σ_1 which determines the level of scatter alignment (rotation and scale/shear) as well as σ_2 which determines the level of alignment of means.

Figure 6 shows that, as $\sigma_1 \rightarrow 0$ and $\sigma_2 \rightarrow 0$, our algorithm converges to the baseline fine-tuning on the combined source and target domains ($S+T$) which yielded 85.9% accuracy for split $sp1$.

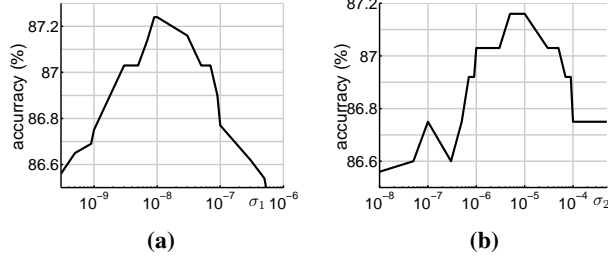


Figure 6: Performance of our second-order alignment loss (So) w.r.t. parameters σ_1 (6a) and σ_2 (6b) on the $\mathcal{A} \rightarrow \mathcal{W}$ domain shift (split $sp1$ is used). Note the logarithmic scale.

Moreover, the overall performance is stable *i.e.*, within $\pm 0.2\%$ accuracy, for a large range of values *e.g.*, $5e-9 \leq \sigma_1 \leq 5e-8$ and $1e-6 \leq \sigma_2 \leq 5e-5$.

B Derivatives of the alignment loss g w.r.t. the feature vectors

Suppose $\Phi = [\phi_1, \dots, \phi_N]$ and $\Phi^* = [\phi_1^*, \dots, \phi_{N^*}^*]$ are some feature vectors of quantity N and N^* , respectively, which are used to evaluate Σ and Σ^* . For $r=2$, we have to first compute the derivative of the covariance matrix Σ w.r.t. $\phi_{m'n'}$. To do so, we proceed by computing derivatives of: i) the autocorrelation matrix in (9) and ii) the outer product of means μ in (10) and (11):

$$\frac{\partial \sum_n \phi_n \phi_n^T}{\partial \phi_{m'n'}} = j_{m'} \phi_{n'}^T + \phi_{n'} j_{m'}^T, \quad (9)$$

$$\frac{\partial \mu \mu^T}{\partial \mu_{m'}} = j_{m'} \mu^T + \mu j_{m'}^T, \quad (10)$$

$$\frac{\partial \mu \mu^T}{\partial \phi_{m'n'}} = \sum_m \frac{\partial \mu \mu^T}{\partial \mu_m} \frac{\partial \mu_m}{\partial \phi_{m'n'}} = \frac{1}{N} (j_{m'} \mu^T + \mu j_{m'}^T), \quad (11)$$

where $j_{m'}$ is a vector of zero entries except for position m' which is equal one. Putting together (9), (10) and (11) yields the derivative of Σ w.r.t. $\phi_{m'n'}$:

$$\frac{\partial \left(\frac{1}{N} \sum_n \phi_n \phi_n^T \right) - \mu \mu^T}{\partial \phi_{m'n'}} = \frac{1}{N} \left(j_{m'} (\phi_{n'} - \mu)^T + (\phi_{n'} - \mu) j_{m'}^T \right). \quad (12)$$

The derivatives of $\|\Sigma - \Sigma^*\|_F^2$ w.r.t. covariance Σ as well as $\phi_{m'n'}$ and $\phi_{m'n'}^*$ are provided below:

$$\frac{\partial \|\Sigma - \Sigma^*\|_F^2}{\partial \Sigma} = 2 (\Sigma - \Sigma^*) \quad (13)$$

$$\begin{aligned} \frac{\partial \|\Sigma - \Sigma^*\|_F^2}{\partial \phi_{m'n'}} &= \sum_{m,n} \frac{\partial \|\Sigma - \Sigma^*\|_F^2}{\partial \Sigma_{mn}} \left(\frac{\partial \Sigma}{\partial \phi_{m'n'}} \right)_{mn} \\ &= \frac{2}{N} \sum_{m,n} (\Sigma - \Sigma^*)_{mn} \left(j_{m'} (\phi_{n'} - \mu)^T + (\phi_{n'} - \mu) j_{m'}^T \right)_{mn} \\ &= \frac{4}{N} (\Sigma_{m',:} - \Sigma_{m',:}^*) (\phi_{n'} - \mu). \end{aligned} \quad (14)$$

The derivatives of $\|\Sigma - \Sigma^*\|_F^2$ w.r.t. Φ and Φ^* are:

$$\frac{\partial \|\Sigma - \Sigma^*\|_F^2}{\partial \Phi} = \frac{4}{N} (\Sigma - \Sigma^*) (\Phi - \mu \mathbf{1}^T), \quad (15)$$

$$\frac{\partial \|\Sigma - \Sigma^*\|_F^2}{\partial \Phi^*} = -\frac{4}{N^*} (\Sigma - \Sigma^*) (\Phi^* - \mu^* \mathbf{1}^T). \quad (16)$$

The derivatives of $\|\mu - \mu^*\|_2^2$ w.r.t. μ , ϕ_n and ϕ_n^* are:

$$\frac{\partial \|\mu - \mu^*\|_2^2}{\partial \mu} = 2 (\mu - \mu^*), \quad (17)$$

$$\frac{\partial \|\mu - \mu^*\|_2^2}{\partial \phi_n} = \frac{2 (\mu - \mu^*)}{N}, \quad \frac{\partial \|\mu - \mu^*\|_2^2}{\partial \phi_n^*} = \frac{2 (\mu - \mu^*)}{N^*}. \quad (18)$$

C Kernelized derivative of the Frobenius norm between tensors w.r.t. the feature vectors

Suppose that some feature vectors $\Phi = [\phi_1, \dots, \phi_N]$ and $\Phi^* = [\phi_1^*, \dots, \phi_{N^*}^*]$ are given in quantities N and N^* and that the Frobenius norm between tensors \mathcal{X} and \mathcal{Y} of order $r \geq 1$ build from Φ and Φ^* is being evaluated. Then, the derivative of Equation (8) w.r.t. feature vector ϕ_{n^\dagger} becomes:

$$\begin{aligned} \frac{\partial \|\mathcal{X} - \mathcal{X}^*\|_F^2}{\partial \phi_{n^\dagger}} &= \frac{1}{N^2} r \sum_{n=1}^N \sum_{n'=1}^N K_{nn'}^{r-1} \frac{\partial K_{nn'}}{\partial \phi_{n^\dagger}} \\ &\quad - \frac{2}{NN^*} r \sum_{n=1}^N \sum_{n'=1}^{N^*} \bar{K}_{nn'}^{r-1} \frac{\partial \bar{K}_{nn'}}{\partial \phi_{n^\dagger}}, \end{aligned} \quad (19)$$

where

$$\begin{aligned} \frac{\partial K_{nn'}}{\partial \phi_{n^\dagger}} &= \frac{\partial \langle \phi_n, \phi_{n'} \rangle}{\partial \phi_{n^\dagger}} - \frac{\partial \langle \mu, \phi_{n'} \rangle}{\partial \phi_{n^\dagger}} - \frac{\partial \langle \phi_n, \mu \rangle}{\partial \phi_{n^\dagger}} + \frac{\partial \langle \mu, \mu \rangle}{\partial \phi_{n^\dagger}} \\ &= \begin{cases} \phi_{n'} & \frac{\phi_{n'}}{N} & (\mu + \frac{\phi_n}{N}) & , \text{ if } n = n^\dagger, n' \neq n^\dagger \\ \phi_n & (\mu + \frac{\phi_{n'}}{N}) & \frac{\phi_n}{N} & , \text{ if } n \neq n^\dagger, n' = n^\dagger \\ 2\phi_n & (\mu + \frac{\phi_{n'}}{N}) & (\mu + \frac{\phi_n}{N}) & , \text{ if } n = n^\dagger, n' = n^\dagger \\ 0 & \frac{\phi_{n'}}{N} & \frac{\phi_n}{N} & , \text{ if } n \neq n^\dagger, n' \neq n^\dagger \end{cases} + \frac{2}{N} \mu \\ \frac{\partial \bar{K}_{nn'}}{\partial \phi_{n^\dagger}} &= \frac{\partial \langle \phi_n, \phi_{n'}^* \rangle}{\partial \phi_{n^\dagger}} - \frac{\partial \langle \mu, \phi_{n'}^* \rangle}{\partial \phi_{n^\dagger}} - \frac{\partial \langle \phi_n, \mu^* \rangle}{\partial \phi_{n^\dagger}} + \frac{\partial \langle \mu, \mu^* \rangle}{\partial \phi_{n^\dagger}} \end{aligned} \quad (20)$$

$$= \begin{cases} \phi_{n'}^* & \mu^* & , \text{ if } n = n^\dagger, n' \neq n^\dagger \\ 0 & 0 & , \text{ if } n \neq n^\dagger, n' = n^\dagger \\ \phi_n^* & -\frac{1}{N} \phi_{n'}^* & - \mu^* + \frac{1}{N} \mu^* & , \text{ if } n = n^\dagger, n' = n^\dagger \\ 0 & 0 & , \text{ if } n \neq n^\dagger, n' \neq n^\dagger \end{cases} \quad (21)$$

Putting together Equations (19), (20) and (21) and setting $q = r - 1$ yields the derivatives w.r.t. matrices Φ and Φ^* :

$$\begin{aligned} \frac{\partial \|\mathcal{X} - \mathcal{X}^*\|_F^2}{\partial \Phi} &= \frac{2}{N^2} r \Phi \left(\mathbf{K}^q - \frac{1}{N} (\mathbb{1}^T \mathbf{K}^q)^T \mathbb{1}^T \right) \\ &\quad + \frac{2r\mu}{N^2} \left(\frac{1}{N} \mathbb{1}^T \mathbf{K}^q \mathbb{1} - \mathbb{1}^T \mathbf{K}^q \right) - \frac{2r\Phi^*}{NN^*} \left(\bar{\mathbf{K}}^q - \frac{1}{N} (\bar{\mathbf{K}}^q \mathbb{1}) \mathbb{1}^T \right) \\ &\quad + \frac{2}{NN^*} r \mu^* \left(\mathbb{1}^T \bar{\mathbf{K}}^q - \frac{1}{N} \mathbb{1}^T \bar{\mathbf{K}}^q \mathbb{1} \right) \end{aligned} \quad (22)$$

and

$$\begin{aligned} \frac{\partial \|\mathcal{X} - \mathcal{X}^*\|_F^2}{\partial \Phi^*} &= \frac{2}{N^{*2}} r \Phi^* \left(\bar{\mathbf{K}}^q - \frac{1}{N^*} (\mathbb{1}^T \bar{\mathbf{K}}^q)^T \mathbb{1}^T \right) \\ &\quad + \frac{2r\mu^*}{N^{*2}} \left(\frac{1}{N^*} \mathbb{1}^T \bar{\mathbf{K}}^q \mathbb{1} - \mathbb{1}^T \bar{\mathbf{K}}^q \right) - \frac{2r\Phi}{NN^*} \left(\bar{\mathbf{K}}^q - \frac{1}{N^*} (\bar{\mathbf{K}}^q \mathbb{1}) \mathbb{1}^T \right) \\ &\quad + \frac{2}{NN^*} r \mu \left(\mathbb{1}^T \bar{\mathbf{K}}^q - \frac{1}{N^*} \mathbb{1}^T \bar{\mathbf{K}}^q \mathbb{1} \right). \end{aligned} \quad (23)$$

References

- [1] Jonathan Baxter, Rich Caruana, Tom Mitchell, Lior Y. Pratt, Daniel L. Silver, and Sebastian Thrun. Learning to learn: Knowledge consolidation and transfer in inductive systems. NIPS Workshop, http://plato.acadiau.ca/courses/comp/dsilver/NIPS95_LTL/transfer_workshop.1995.html, 1995. Accessed: 30-10-2016.

- [2] Lin Chen, Wen Li, and Dong Xu. Recognizing rgb images by learning from rgb-d data. *CVPR*, 2014.
- [3] Sumit Chopra, Suhril Balakrishnan, and Raghuraman Gopalan. Dlid: Deep learning for domain adaptation by interpolating between domains. *ICML Workshop*, 2013.
- [4] K. Crammer, M. Kearns, and J. Wortman. Learning from multiple sources. *JMLR*, 9:1757–1774, 2008.
- [5] Hal Daumé, III, Abhishek Kumar, and Avishek Saha. Frustratingly easy semi-supervised domain adaptation. *Proceedings of the 2010 Workshop on Domain Adaptation for Natural Language Processing*, pages 53–59, 2010.
- [6] Jeff Donahue, Yangqing Jia, Oriol Vinyals, Judy Hoffman, Ning Zhang, Eric Tzeng, and Trevor Darrell. Decaf: A deep convolutional activation feature for generic visual recognition. *ICML*, 2014.
- [7] Mathias Eitz, James Hays, and Marc Alexa. How do humans sketch objects? *ACM Trans. Graph. (Proc. SIGGRAPH)*, 31(4):44:1–44:10, 2012.
- [8] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2007 (VOC2007) Results. <http://pascallin.ecs.soton.ac.uk/challenges/VOC, 2007>.
- [9] Kunihiro Fukushima. Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position. *Biological Cybernetics*, 36(4):193–202, 1980. doi: 10.1007/BF00344251.
- [10] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *JMLR*, 17(1):2096–2030, 2016. ISSN 1532-4435.
- [11] M. Ghifary, W. B. Kleijn, and M. Zhang. Domain adaptive neural networks for object recognition. *CoRR*, abs/1409.6041, 2014.
- [12] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. *CVPR*, pages 580–587, 2014. doi: 10.1109/CVPR.2014.81. URL <http://dx.doi.org/10.1109/CVPR.2014.81>.
- [13] G. Griffin, A. Holub, and P. Perona. Caltech-256 object category dataset. Technical Report 7694, California Institute of Technology, 2007. URL <http://authors.library.caltech.edu/7694>.
- [14] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.
- [15] J. J. Hopfield. Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the National Academy of Sciences*, 79:2554–2558, 1982.
- [16] Nathan Intrator and Shimon Edelman. Making a low-dimensional representation suitable for diverse tasks. *Connection Sci.*, 8(2):205–223, 1996.
- [17] Tae-Kyun Kim, Kwan-Yee Kenneth Wong, and R. Cipolla. Tensor canonical correlation analysis for action classification. *CVPR*, 2007.
- [18] P. Koniusz and A. Cherian. Sparse coding for third-order super-symmetric tensor descriptors with application to texture recognition. *CVPR*, 2016.
- [19] P. Koniusz, A. Cherian, and F. Porikli. Tensor representations via kernel linearization for action recognition from 3D skeletons. *ECCV*, 2016.
- [20] P. Koniusz, F. Yan, P. Gosselin, and K. Mikolajczyk. Higher-order occurrence pooling for bags-of-words: Visual concept detection. *PAMI*, 2016.
- [21] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. ImageNet classification with deep convolutional neural networks. *NIPS*, pages 1106–1114, 2012.
- [22] Ilya Kuzborskij, Fabio Maria Carlucci, and Barbara Caputo. When naïve bayes nearest neighbors meet convolutional neural networks. *CVPR*, 2016.
- [23] R.; Perona L. Fei-Fei; Fergus. One-shot learning of object categories. *TPAMI*, 28:594–611, April 2006.
- [24] Kevin Lai, Liefeng Bo, Xiaofeng Ren, and Dieter Fox. A large-scale hierarchical multi-view rgb-d object dataset. *ICRA*, pages 1817–1824, 2011. URL <http://dblp.uni-trier.de/db/conf/icra/icra2011.html#LaiBRF11>.
- [25] L. De Lathauwer, B. De Moor, and J. Vandewalle. A multilinear singular value decomposition. *SIAM J. Matrix Analysis and Applications*, 21:1253–1278, 2000.
- [26] W. Li, T. Tommasi, F. Orabona, D. Vázquez, M. López, J. Xu, and H. Larochelle. Task-cv: Transferring and adapting source knowledge in computer vision. *ECCV Workshop*, <http://adas.cvc.uab.es/task-cv2016>, 2016. Accessed: 22-11-2016.
- [27] H. Lu, K. N. Plataniotis, and A. N. Venetsanopoulos. A survey of multilinear subspace learning for tensor data. *Pattern Recognition*, 44(7):1540–1551, 2011.

- [28] Saeid Motiian and Gianfranco Doretto. Information bottleneck domain adaptation with privileged information for visual recognition. *ECCV*, 2016.
- [29] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet large scale visual recognition challenge. *IJCV*, 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.
- [30] Kate Saenko, Brian Kulis, Mario Fritz, and Trevor Darrell. Adapting visual category models to new domains. *ECCV*, pages 213–226, 2010. URL <http://dl.acm.org/citation.cfm?id=1888089.1888106>.
- [31] Pierre Sermanet, David Eigen, Xiang Zhang, Michael Mathieu, Rob Fergus, and Yann Lecun. Overfeat: Integrated recognition, localization and detection using convolutional networks. *ICLR*, 2014. URL <http://arxiv.org/abs/1312.6229>.
- [32] A. Shashua and T. Hazan. Non-negative tensor factorization with applications to statistics and computer vision. *ICML*, 2005.
- [33] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *ICLR*, abs/1409.1556, 2015.
- [34] Baochen Sun, Jiashi Feng, and Kate Saenko. Return of frustratingly easy domain adaptation. *CoRR*, abs/1511.05547, 2015. URL <http://arxiv.org/abs/1511.05547>.
- [35] Sebastian Thrun. Is learning the n-th thing any easier than learning the first? *NIPS*, pages 640–646, 1996.
- [36] Tatiana Tommasi, Francesco Orabona, and Barbara Caputo. Safety in numbers: Learning categories from few examples with multi model knowledge transfer. *CVPR*, pages 3081–3088, 2010. doi: 10.1109/CVPR.2010.5540064.
- [37] Tatiana Tommasi, Martina Lanzi, Paolo Russo, and Barbara Caputo. Learning the roots of visual domain shift. *ECCV Workshop*, 2016.
- [38] E. Tzeng, J. Hoffman, T. Darrell, and K. Saenko. Simultaneous deep transfer across domains and tasks. *ICCV*, pages 4068–4076, 2015.
- [39] Vladimir Vapnik and Akshay Vashist. A new learning paradigm: Learning using privileged information. *Neural Networks*, 22(56):544–557, 2009. doi: <http://dx.doi.org/10.1016/j.neunet.2009.06.042>.
- [40] M. A. Vasilescu and D. Terzopoulos. Multilinear analysis of image ensembles: Tensorfaces. *ECCV*, 2002.
- [41] M. A. Vasilescu and D. Terzopoulos. Tensortextures: multilinear image-based rendering. *ACM Transactions on Graphics*, 23(3):336–342, 2004.
- [42] Yu-Xiong Wang and Martial Hebert. Learning to learn: Model regression networks for easy small sample learning. *ECCV*, 2016.
- [43] J. West, D. Venture, and S. Warnick. Spring research presentation: A theoretical foundation for inductive transfer. *Brigham Young University, College of Physical and Mathematical Sciences*, 2007.
- [44] R. S. Woodworth and E. L. Thorndike. The influence of improvement in one mental function upon the efficiency of other functions. *Psychological Review (1)*, 8(3):247–261, 1901. doi: 10.1037/h0074898.
- [45] Yi-Ren Yeh, Chun-Hao Huang, and Yu-Chiang Frank Wang. Heterogeneous domain adaptation and classification by exploiting the correlation subspace. *Transactions on Image Processing*, 23(5), 2014.